



AFRL-RI-RS-TR-2010-104

**COMBINING FACTS AND EXPERT OPINION IN ANALYTICAL MODELS VIA  
LOGICAL AND PROBABILISTIC REASONING**

---

University of South Carolina Research Foundation

*April 2010*

FINAL TECHNICAL REPORT

*APPROVED FOR PUBLIC RELEASE; DISTRIBUTION UNLIMITED.*

STINFO COPY

**AIR FORCE RESEARCH LABORATORY  
INFORMATION DIRECTORATE**

## NOTICE AND SIGNATURE PAGE

Using Government drawings, specifications, or other data included in this document for any purpose other than Government procurement does not in any way obligate the U.S. Government. The fact that the Government formulated or supplied the drawings, specifications, or other data does not license the holder or any other person or corporation; or convey any rights or permission to manufacture, use, or sell any patented invention that may relate to them.

This report was cleared for public release by the 88<sup>th</sup> ABW, Wright-Patterson AFB Public Affairs Office and is available to the general public, including foreign nationals. Copies may be obtained from the Defense Technical Information Center (DTIC) (<http://www.dtic.mil>).

AFRL-RI-RS-TR-2010-104 HAS BEEN REVIEWED AND IS APPROVED FOR PUBLICATION IN ACCORDANCE WITH ASSIGNED DISTRIBUTION STATEMENT.

FOR THE DIRECTOR:

/s/  
NANCY A. ROBERTS  
Work Unit Manager

/s/  
JOSEPH CAMERA, Chief  
Information & Intelligence Exploitation Division  
Information Directorate

This report is published in the interest of scientific and technical information exchange, and its publication does not constitute the Government's approval or disapproval of its ideas or findings.

**REPORT DOCUMENTATION PAGE***Form Approved*  
**OMB No. 0704-0188**

Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden to Washington Headquarters Service, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188) Washington, DC 20503.

**PLEASE DO NOT RETURN YOUR FORM TO THE ABOVE ADDRESS.****1. REPORT DATE (DD-MM-YYYY)**

APRIL 2010

**2. REPORT TYPE**

Final

**3. DATES COVERED (From - To)**

August 2006 – December 2009

**4. TITLE AND SUBTITLE**COMBINING FACTS AND EXPERT OPINION IN ANALYTICAL  
MODELS VIA LOGICAL AND PROBABILISTIC REASONING**5a. CONTRACT NUMBER**

FA8750-06-C-0194

**5b. GRANT NUMBER**

N/A

**5c. PROGRAM ELEMENT NUMBER**

31011G

**6. AUTHOR(S)**

Marco Valtorta, Michael Huhns, and John Byrnes

**5d. PROJECT NUMBER**

CASE

**5e. TASK NUMBER**

00

**5f. WORK UNIT NUMBER**

06

**7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)**University of South Carolina Research Foundation  
901 Sumter Street, Suite 511  
Columbia, SC 29208-0001**8. PERFORMING ORGANIZATION  
REPORT NUMBER**

N/A

**9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)**AFRL/RIEH  
525 Brooks Road  
Rome NY 13441-4505**10. SPONSOR/MONITOR'S ACRONYM(S)**

N/A

**11. SPONSORING/MONITORING  
AGENCY REPORT NUMBER**  
AFRL-RI-RS-TR-2010-104**12. DISTRIBUTION AVAILABILITY STATEMENT**

Approved for Public Release; Distribution Unlimited. PA# 88ABW-2010-2087

Date Cleared: 20-April-2010

**13. SUPPLEMENTARY NOTES****14. ABSTRACT**

This report describes work on an integrated system that can assist analysts in exploring hypotheses using Bayesian analysis of evidence from a variety of sources. The hypothesis exploration is aided by an ontology that represents domain knowledge, events, and causality for Bayesian reasoning, as well as models of information sources for evidential reasoning. We are validating the approach via a tool, Magellan, that uses both Bayesian models and logical models for an analyst's prior and tacit knowledge about how evidence can be used to evaluate hypotheses. We also describe how we combine logic information, in the form of proofs provided by the natural deduction system SILK(Semantic Inferencing on Large Knowledge) and probabilistic information, represented by Bayesian networks, in the BRUSE(Bayesian Reasoning Using Soft Evidence) system.

**15. SUBJECT TERMS**

Bayesian and logical models, Bayesian reasoning, hypothesis analysis, ontology management, uncertainty, probabilistic models

**16. SECURITY CLASSIFICATION OF:****a. REPORT**  
U**b. ABSTRACT**  
U**c. THIS PAGE**  
U**17. LIMITATION OF  
ABSTRACT**

UU

**18. NUMBER  
OF PAGES**

34

**19a. NAME OF RESPONSIBLE PERSON**

Nancy A. Roberts

**19b. TELEPHONE NUMBER (Include area code)**

N/A

## **Acknowledgments**

This work was funded in part by the Disruptive Technology Office Collaboration and Analyst System Effectiveness (CASE) Program, contract FA8750-06-C-0194 issued by Air Force Research Laboratory (AFRL). The views and conclusions are those of the authors, not of the US Government or its agencies.

# Table of Contents

1.0	Introduction.....	1
2.0	Bayesian Reasoning for Evidence Management.....	2
2.1	Recognizing and Representing Situations.....	2
2.2	Capturing the Terminology and Prior Knowledge for a New Domain.....	4
2.3	Situation Fragments Represented by Logical Models.....	5
2.4	An Extended Example.....	7
2.5	Causality.....	14
2.6	Evidence .....	16
3.0	Use of Tripartite Ontology for Intelligence Analysis .....	19
4.0	Evaluation .....	22
5.0	Conclusion .....	23
5.1.1	The Relationship between Proofs and Bayesian Networks .....	23
5.1.2	Proofs and Theories .....	24
5.1.3	First-order Theories .....	24
5.1.4	Semantics .....	24
6.0	References.....	26
7.0	List of Acronyms .....	27

## List of Figures

Figure 1: A commonly occurring part of a situation for a suspicious bank transfer of money, represented as an uninstantiated Bayesian network.....	3
Figure 2: Fragments (templates) are merged based on the evidence that instantiates them .....	4
Figure 3: An ontology for intelligence analysts has three related parts, corresponding to (1) the world of causality and hypothetical events needed for Bayesian reasoning, (2) the real world of things needed to model situations, and (3) the world of information and information sources needed for evidence management. ....	5
Figure 4: The BALER framework for integrating logical models with probabilistic models, with an ontology developed in Protégé providing a consistent vocabulary for all domain concepts .....	6
Figure 5: An example illustrating the need for both Bayesian and logical reasoning .....	7
Figure 6: An example illustrating the need for both Bayesian and logical reasoning .....	8
Figure 7: This proof tree makes contexts explicit. $\Gamma$ stands for the set of assumptions $\{C \rightarrow B, T \rightarrow B, C \vee T\}$ , and $\Gamma, A$ stands for $\Gamma \cup \{A\}$ . ....	8
Figure 8: B logically follows from the axioms in the brown liquids domain .....	8
Figure 9: A probabilistic causal model that relates work deadlines to coffee and tea in my cup	10
Figure 10: A model composed from logical and probabilistic components .....	10
Figure 11: Probability update in the model of the previous figure .....	11
Figure 12: Coffee and tea may not be together in the cup .....	12
Figure 13: Probability update in the model of the previous figure .....	12
Figure 14: Soft evidential update with a 10% exception rate for the constraint that only one drink may be present in the cup .....	13
Figure 15: The BALER software process flow, which is supported by the tripartite ontology of real world concepts, events, and information sources .....	14
Figure 16: Protégé is used to enter the ontology concepts that form the basis for representing situations and evidence .....	15
Figure 17: Evidence consists of a set of findings, which can be of three different types.....	16
Figure 18: Magellan’s extended ACH interface is integrated with the ontology of events through pull-down menus.....	17

Figure 19: A Bayesian network fragment constructed automatically from an ACH matrix. The conditional probabilities needed for Bayesian reasoning are derived from the user-entered values in the matrix indicating whether or not a finding is consistent with an analyst's hypothesis.....	18
Figure 20: Magellan architecture for Bayesian Reasoning used to explore an analyst's hypotheses.....	19
Figure 21: A small portion of the tripartite ontology indicating how an item of evidence would be classified and used to instantiate one or more fragments.....	20
Figure 22: The Value of Information calculation algorithm used in Magellan. ....	21
Figure 23: The Magellan interface showing an evidence message, the ontology concepts it contains, the fragments that it instantiates, composed into a situation, and the posterior probability for an hypothesis about the situation.....	21

## **1.0 Introduction**

Much of the extensive work on ontologies to date has focused on modeling and representing the world of objects. The ontologies needed for our research supporting the management of hypotheses and evidence for analysts, however, must additionally model events and causality. Less work has been done on this aspect of ontologies. In this paper we show how concepts from a causal ontology can be used directly as variables in Bayesian networks and how the attributes of the causal concepts can be used in matching evidence to the variables. Moreover, subclass relationships in the ontology enable the extension of Bayesian reasoning over types.



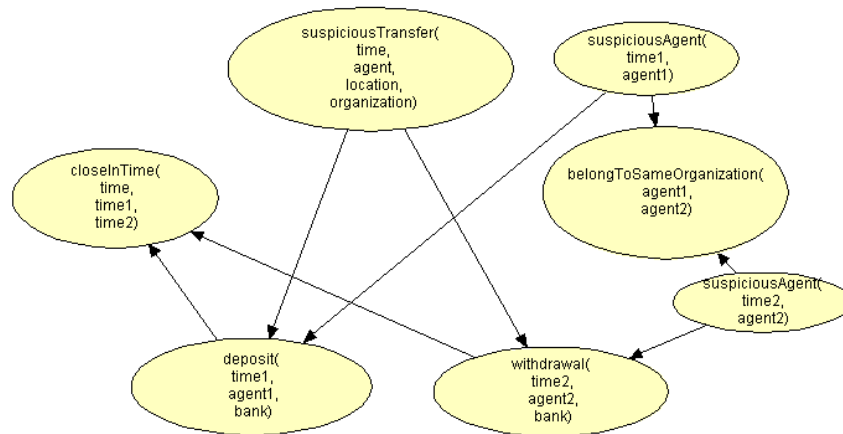
## 2.0 Bayesian Reasoning for Evidence Management

There are numerous real-world situations about which an analyst might wish to hypothesize and investigate, but it would be impractical to encode all of them explicitly in a support system for analysts. Instead, our approach is to represent fragments of situations and provide a mechanism for combining them into a wide variety of more complete ones [1,4]. The combination occurs dynamically as evidence about a situation becomes available or as an analyst revises or enters new hypotheses. A situation fragment is represented as a Bayesian network with nodes for hypotheses, events, and evidence, and links for relating them. Our ability to combine the fragments into more complete situation models is dependent on having a consistent terminology in which the fragments are described. The focus of our work has been on (1) defining and representing the terminology, including terms of a domain and terms for evidence in that domain, (2) capturing new fragments from a variety of sources, and (3) incorporating the terminology and BN fragments into an integrated end-to-end tool, Magellan.

### 2.1 Recognizing and Representing Situations

Our objective is to be able to model and reason probabilistically about a wide variety of situations that might be of interest to analysts. Unfortunately, there are too many of these to encode *a priori* within a system for analysts and they are too complex for most analysts to encode *a posteriori* complex real-world situations formally is unrealistic. Instead, our approach is to represent small, common aspects of situations generically, and then provide a means to combine them dynamically into representations for real-world situations. We term the small generic situation aspect a *fragment*, and choose a first-order representation for it.

An example situation aspect that we might represent as a fragment would be a “suspicious transfer of money,” with variables corresponding to banks, organizations, deposits, withdrawals, and the transferring agent. The fragment would be instantiated when evidence matched the variables, e.g., “a church attended by Syrians in Detroit deposited funds into a Michigan bank and the funds were transferred to a bank in Cairo.” More precisely, each variable (node) in a fragment has a set of identifying attributes and their collective instantiated values specify a particular instance of a random variable. Because the evidence might be uncertain, there would be probabilities associated with the instantiated fragment, and we would treat the instantiated fragment as a Bayesian network. This is shown in Figure 1. Note that the probability distribution described in the Bayesian network is a joint distribution on the nodes only, not on the nodes and the attributes.



**Figure 1. A commonly occurring part of a situation for a suspicious bank transfer of money, represented as an uninstantiated Bayesian network**

An advantage of using fragments of situations instead of more complete situations is that many more situations can be represented efficiently. More precisely,  $N$  fragments can potentially be combined in  $N!$  ways to represent  $N!$  situations. The combining is guided by available evidence. For example, three other situations that we might represent as fragments are “purchases of weapons,” “influencing an election,” and “bribing a politician.” If evidence matched one of these, and the resulting instantiated fragment had one or more variables in common with the money fragment, then we would merge the fragments at the point of the common variables to produce a representation of a more complete situation, such as “transferring money to influence an election.” Note that fragments can be merged only if the attributes of their common variables unify. Also note that it is not necessary for the fragments to have any variables in common in order to merge them and represent larger situations. As a result, the fragments could represent situations such as “bribing a politician to influence an election” and “purchasing weapons to influence an election.” Further, because each fragment could be instantiated multiple times, we could represent several different money transfers being used to purchase weapons. Using Magellan, Bayesian reasoning would then be performed on whichever complex situation representation resulted from instantiating fragments with the available evidence and integrating those fragments. The overall process for merging instantiated fragments and reasoning over them is shown in Figure 2.

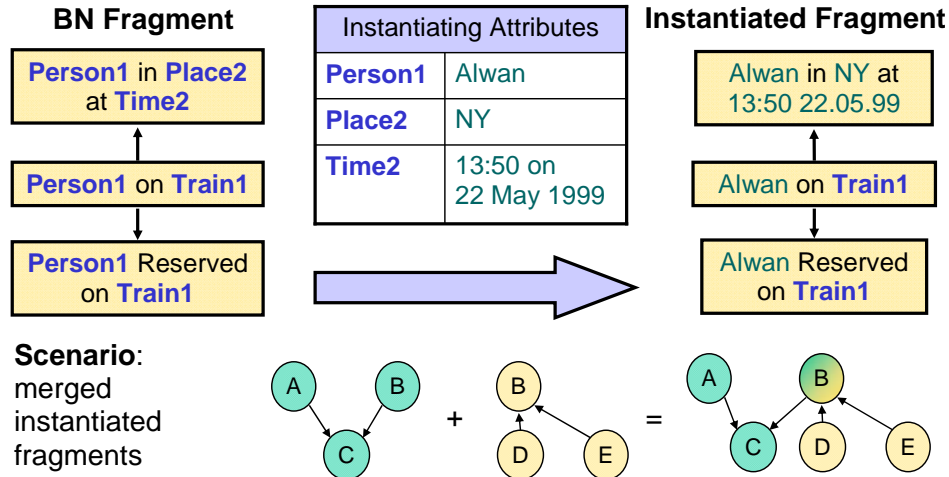
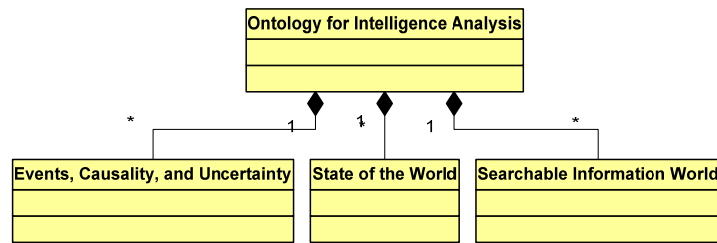


Figure 2. Fragments (templates) are merged based on the evidence that instantiates them

## 2.2 Capturing the Terminology and Prior Knowledge for a New Domain

A key activity of an intelligence analyst is to distinguish among competing hypotheses, determine the likelihood of their occurrence, and reduce the uncertainty in the outcomes of the hypotheses, upon which decision makers will then base their decisions. Hypothesis outcomes<sup>1</sup> are related to observable evidence via direct or indirect causal relations, and therefore ontological support for analysts should involve cause-and-effect. This is best supported by an ontology emphasizing events and their causal relationships, along with a hypothetical world of possible events, actions, and causes. However, causal relationships must be interpreted in the context of the state of the real world—primarily consisting of objects and their physical properties—which can be represented in a conventional ontology, such as those that are part of SUMO. The evidence for reasoning about hypotheses can come from a variety of sources, and the acquisition of evidence and events from these sources must also be represented, constituting a third kind of ontological representation describing the information sources. Figure 3 depicts the three ontological models we use for (1) modeling situations and relating them to (2) background knowledge about the state of the world, and (3) acquiring evidence, all of which enables an assessment of the likelihood of the situations using Bayesian reasoning.

<sup>1</sup> In our ontology, an outcome is thus an important and necessary property (“slot” in Protégé) for hypotheses and, indeed, for any concept that may be in a causal relationship. The relationship is a link in a Bayesian network.



**Figure 3. An ontology for intelligence analysts has three related parts, corresponding to (1) the world of causality and hypothetical events needed for Bayesian reasoning, (2) the real world of things needed to model situations, and (3) the world of information and information sources needed for evidence management.**

A situation might represent an analyst’s query or, more generally, provide context and support for a hypothesis. A situation would be comprised of one or more items of interest and each such item of interest has information provided by several information sources. An item of interest may be specialized to Person, Organization, Event, or Place, and of particular interest would be items relating events involving people at significant places. Information sources can be maps, images, reports video, audio, email, websites, and database records. Typically, an item of interest would have many information sources describing aspects of that item, for example a meeting held by members of a suspected terrorist organization might be described by audio, video, and email surveillance or reports by insiders. Our tool, Magellan, uses Protégé (see Figure 4) for capturing the ontologies, RDF for representing the terminology, XMLBIF for representing the causal relationships, and RDF and SPARQL (future) for requesting evidence from information sources. It also makes use of logical, non-probabilistic models, as shown in Figure 4 and described next.

### **2.3 Situation Fragments Represented by Logical Models**

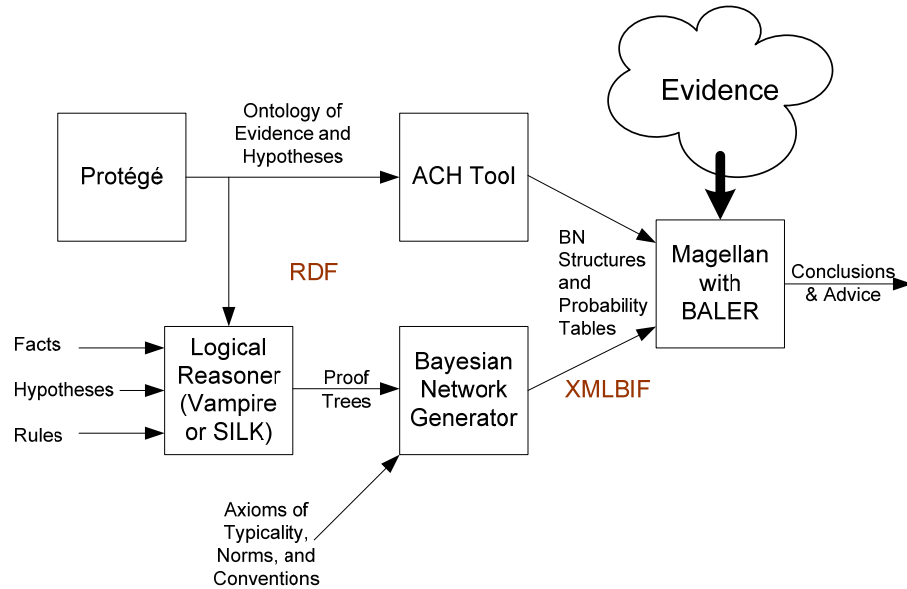
Our objective is to produce models of systems and situations that will be sufficiently accurate that they can be used—where appropriate—to predict future states, to understand operations, to illuminate the factors relevant to decisions, and to control behaviors. We have realized that some knowledge is more easily and naturally represented in the form of statements in a logic language and some is more naturally represented in a Bayesian-network formalism. For example, logic is best for expressing

- Class-subclass statements, such as “C4 is an explosive”
- Part-whole statements, such as “triggers are part of IEDs”
- Definitional statements, such as “triangles have three sides”
- Temporal statements, such as “3:00 p.m. occurs before 4:00 p.m.”
- Spatial statements, such as “Irbil is located in Kurdistan”

Other knowledge is probabilistic, such as

- “Terrorist cell X planned the bombing”
- “Suspect Y met with cell leader Z in Syria last March”

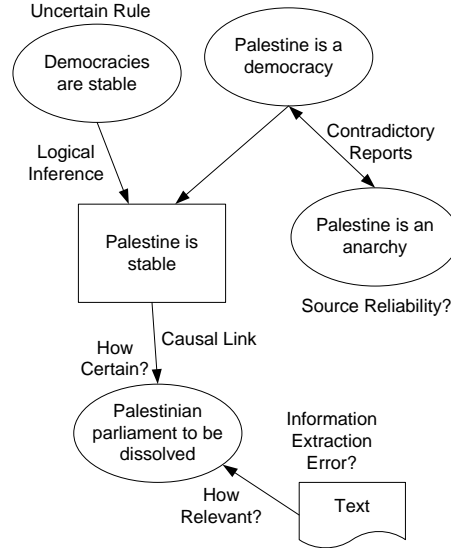
Our objective has been to take advantage of the strengths of each formalism while combining them into a single coherent system.



**Figure 4. The BALER framework for integrating logical models with probabilistic models, with an ontology developed in Protégé providing a consistent vocabulary for all domain concepts**

An example of the situations that can be represented by such an integrated system is shown in Figure 5. This system would help analysts confront problems of *credibility*, *relevance*, *contradictory evidence*, and *pervasive uncertainty*, using

- A unique combination of the power of logical and probabilistic reasoning
- Numerical analysis of competing hypotheses
- Automated linking of relevant evidence
- Automated propagation of uncertainty values: good arguments from uncertain data still add strength to a conclusion
- Robust reasoning over contradictory information allows analysts to exploit maximal amounts of information
- A provision for analysts to enter their own knowledge directly, allowing the system to learn from its users
- The use of probabilities to quantify belief in hypotheses to support optimal decision making according to the principle of maximum expected utility.



**Figure 5. An example illustrating the need for both Bayesian and logical reasoning**

Formal logical tools are able to provide some amount of reasoning support for information analysis, but are unable to represent uncertainty. Bayesian network tools represent probabilistic and causal information, but in the worst case they scale as poorly as some formal logical systems and require specialized expertise to use effectively. The framework (BALER) we have developed for intelligence reasoning incorporates the advantages of both Bayesian and logical systems [7]. The framework includes a formal mechanism for the conversion of automatically generated natural deduction proof trees into Bayesian networks. This is indicated by the information flow shown in Figure 4. We have proven that the merging of such networks with domain-specific causal models forms a consistent Bayesian network with correct values for the formulas derived in the proof. In particular, we show that when the premises of a proof are true, hard evidential update (see Section 2.5) forces the conclusions of the proof to be true with probability one, regardless of any dependencies and prior probability values assumed for the causal model.

## 2.4 An Extended Example

We provide an extended example of using the integrated logical and probabilistic reasoning system. Since the propositional theory that formalizes the example includes at least one non-Horn clause, i.e., at least one clause that includes two non-negative literals, the theory cannot be handled correctly by Prolog or by forward chaining rule-based systems such as JESS or CLIPS. The example formalizes the following story: my cup contains either coffee (C) or tea (T). Coffee is a brown liquid (B). Tea is a brown liquid. Thus it can be concluded that my cup contains a brown liquid. The axioms in the knowledge base that formalizes the story are:

English Assertion	Logical Representation
My cup contains either coffee or tea	$C \vee T$
Coffee is a brown liquid	$C \rightarrow B$
Tea is a brown liquid	$T \rightarrow B$

We want to show  $B$ . Note that the theory allows for both tea and coffee to be in my cup at the same time. A natural deduction proof for  $B$  is given in Figure 6. The proof consists of three steps: two  $\rightarrow$ -elimination steps and one  $\vee$ -elimination step. The  $\rightarrow$ -elimination steps require one assumption each, namely  $C$  and  $T$ . The  $\vee$ -elimination step, which corresponds to a case analysis step, discharges the assumptions made in the  $\rightarrow$ -elimination steps.

$$\frac{C \vee T \quad \frac{C \rightarrow B \quad C^{(1)}}{B} \quad \frac{T \rightarrow B \quad T^{(2)}}{B}}{B} (1,2)$$

Figure 6. An example illustrating the need for both Bayesian and logical reasoning

An issue that had to be resolved is that of representing the proof in a convenient machine-readable form. Due to the prevalence of XML, we decided to use a variation of the XML format used in the Vampire theorem prover [17]. Since Vampire is a resolution theorem prover, while we use natural deduction, we modified the schema by allowing for an explicit representation of the rule used and of the *context*, defined as the set of assumptions, used in a proof step.

$$\frac{\Gamma \vdash C \vee T \quad \frac{\Gamma \vdash C \rightarrow B \quad \Gamma, C \vdash C}{\Gamma, C \vdash B} \rightarrow\text{-Elim} \quad \frac{\Gamma \vdash T \rightarrow B \quad \Gamma, T \vdash T}{\Gamma, T \vdash B} \rightarrow\text{-Elim}}{\Gamma \vdash B} \vee\text{-Elim}$$

Figure 7. This proof tree makes contexts explicit.  $\Gamma$  stands for the set of assumptions  $\{C \rightarrow B, T \rightarrow B, C \vee T\}$ , and  $\Gamma, A$  stands for  $\Gamma \cup \{A\}$ .

Figure 7 presents the same proof as Figure 6, but in a way that emphasizes the contexts used. The proof in Figure 7, which includes the  $\vdash$  symbol, will remind some readers of the sequent calculus. However, it is directly a natural deduction proof with the exact same structure as the proof in Figure 6; which only uses a different syntax to denote active assumptions.

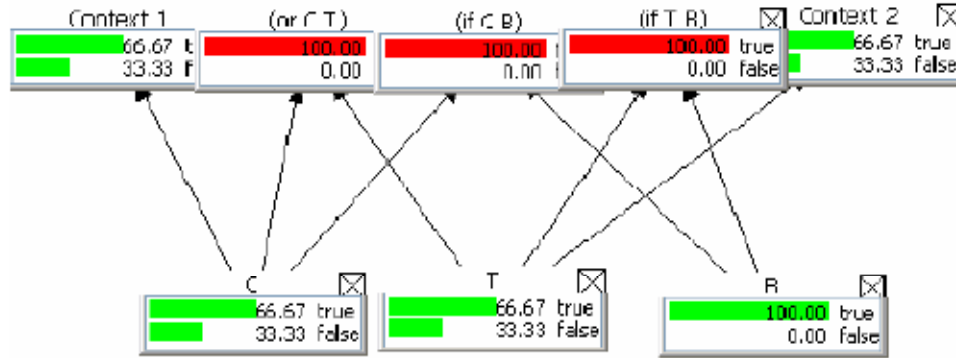


Figure 8.  $B$  logically follows from the axioms in the brown liquids domain

The natural deduction proof is converted to a Bayesian network in the following way. Each non-atomic formula used in the proof, is the child of its component subformulas, with a conditional probability table (CPT) that encodes the main connective introduced or eliminated. For example, in Figure 8, the (nodes corresponding to the) atomic formulas  $C$  and  $T$  are parents of the (node corresponding to the) formula  $(\text{or } C \text{ } T)$ , and the CPT for the family of those three nodes,  $P((\text{or } C \text{ } T) \mid C, T)$  is an OR table. The Bayesian network also represents the nonempty contexts (sets of assumptions) used in the proof. For example, formula  $C$  is the context for the first step of the proof, namely the implication elimination with premises  $C$  and  $(\text{if } C \text{ } B)$  and conclusion  $B$ . Accordingly, the node corresponding to formula  $C$  is a parent of the node  $\text{Context1}$  in the Bayesian network. In the CPT for a context node, the context is true if and only if all of its parents are true.

The construction algorithm just outlined ensures that any possible (i.e., non-zero probability) configuration (i.e., assignment of truth or false values) of the variables in the Bayesian network that correspond to formulas is a true interpretation (a model) of the formulas that appear in the steps of the proof and that no other assignments have positive probability, when the value true is entered as evidence for the (nodes corresponding to the) formulas of the theory. Figure 8 illustrates this, where it is shown that the only state of positive probability of the  $B$  variable is the one in which  $B$  is true, when evidence is entered for  $(\text{if } T \text{ } B)$ ,  $(\text{if } C \text{ } B)$ , and  $(\text{or } C \text{ } T)$ . (Evidence entered is indicated by red bars in a color version of the figure.) Moreover, for a particular set of contexts, the possible configurations are models of the assumptions in the contexts and of the formulas.

Now, imagine that we have probabilistic information relating some of the variables in our domain of interest. In particular, following our example, imagine a probabilistic causal model is available that relates the presence of tea or coffee in my cup to the amount of work I need to get done before the end of the workday, as described in Figure 9. We can now compose the logically derived model of Figure 8 and the probabilistic causal model of Figure 9 into a single model using the Bayesian network fragment composition algorithm described in [1] and obtain the combined model of Figure 10. The combined model is a Bayesian network and can be subjected to processing as any such network. The most important kind of processing is to compute the posterior probability of each variable in the network given a set of findings (i.e., evidence). For example, we may be interested in the probability of a deadline given that we observe coffee in my cup and that all axioms hold (to meet deadlines we work late and consume coffee to stay awake). The posterior probabilities, computed using the commercial Bayesian network shell Hugin ([www.hugin.com](http://www.hugin.com)), are shown in Figure 11, where we observe a roughly 64% probability of my working on a deadline, which happens to be quite a bit higher than the baseline in the model.



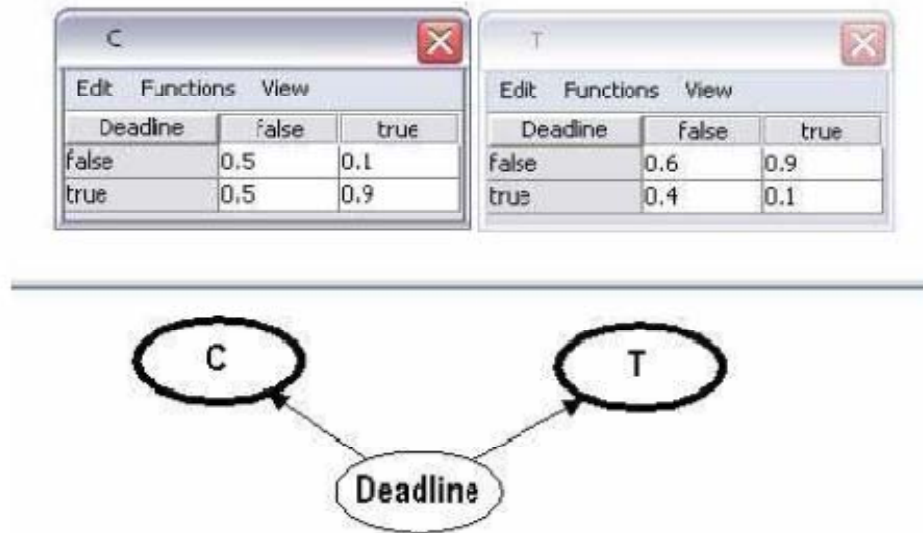


Figure 9. A probabilistic causal model that relates work deadlines to coffee and tea in my cup

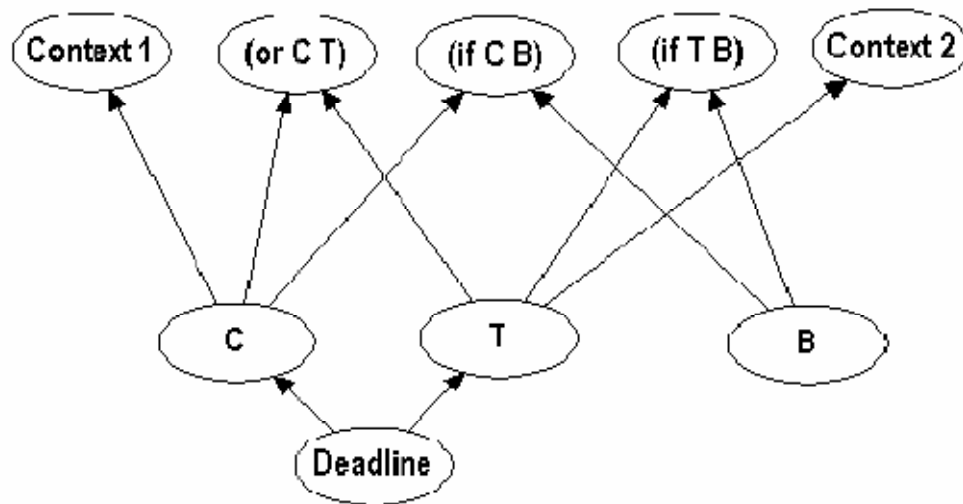
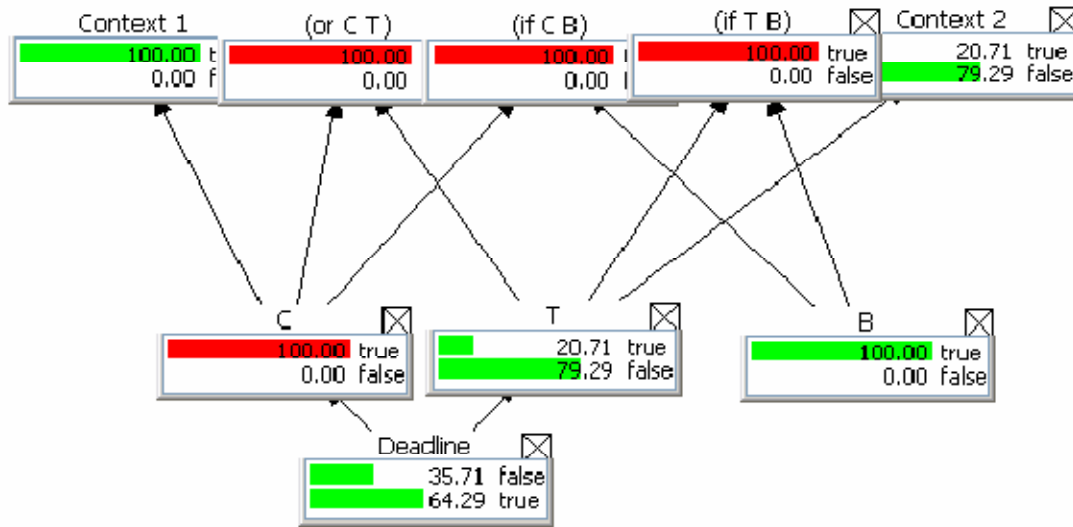


Figure 10. A model composed from logical and probabilistic components



**Figure 11. Probability update in the model of the previous figure**

We also want to allow probability update in the presence of information about the probability of the formulas in the network. For this purpose, we use BRUSE, a refinement of BC-Hugin, a shell that allows the specification of evidence in the form of a set of findings, where each finding is a marginal probability on a variable in the network. In this way, one can specify the probability of a formula holding in the network. BRUSE computes rather efficiently the posterior distribution of the variables in the network with the following properties: the distribution is the closest one (according to cross-entropy) to the original one for which (1) all findings hold, and (2) all d-separation conditions hold. Suppose that, in our example, we add the information that my cup may not contain both coffee and tea at the same time. For simplicity, rather than expressing this constraint as a logical axiom ( $\neg (C \& T)$ ) and converting it into a Bayesian network, we encode this information directly in the Bayesian network, as shown in Figure 12. Comparing Figure 13 with Figure 11 shows that the probability of working under a deadline given that there is coffee in my cup and all axioms hold is about 10% higher than before adding the constraint. Moreover, assume that there are exceptions to this rule. Figure 14 shows the result of running BRUSE on the network, with the exceptions to the rule quantified at 10%. Figure 15 shows the overall logical and probabilistic reasoning process.

OneOnly				
Edit Functions View				
T	true		false	
C	true	false	true	false
false	1	0	0	0.5
true	0	1	1	0.5

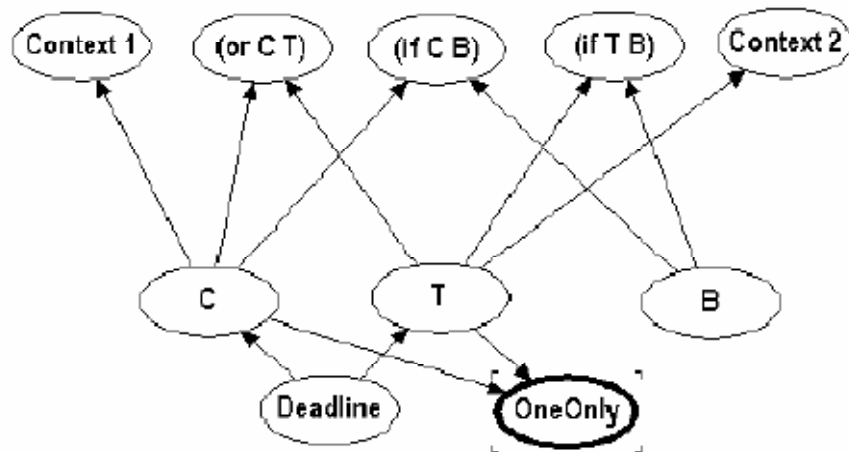


Figure 12. Coffee and tea may not be together in the cup

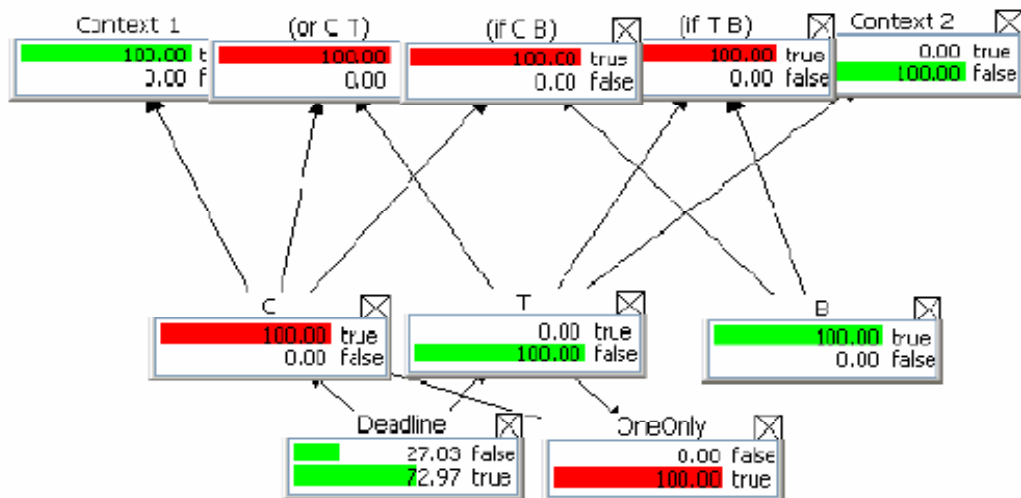


Figure 13. Probability update in the model of the previous figure

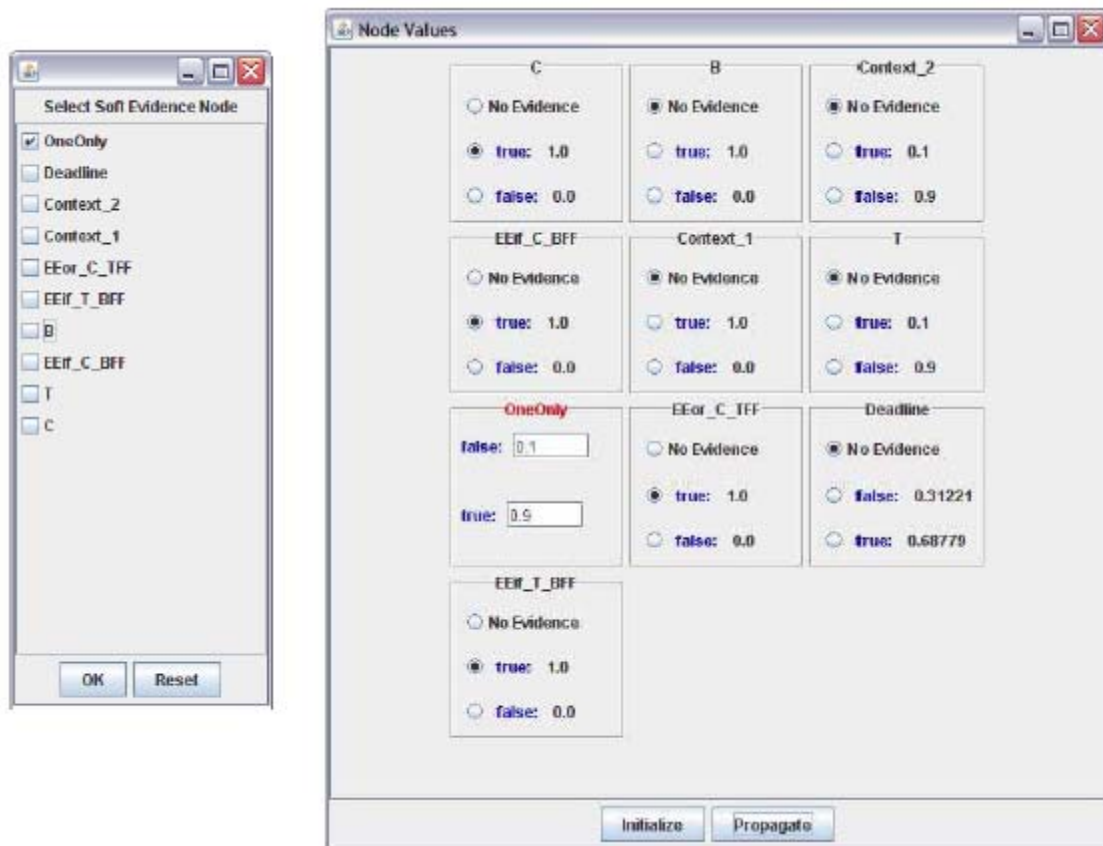


Figure 14. Soft evidential update with a 10% exception rate for the constraint that only one drink may be present in the cup

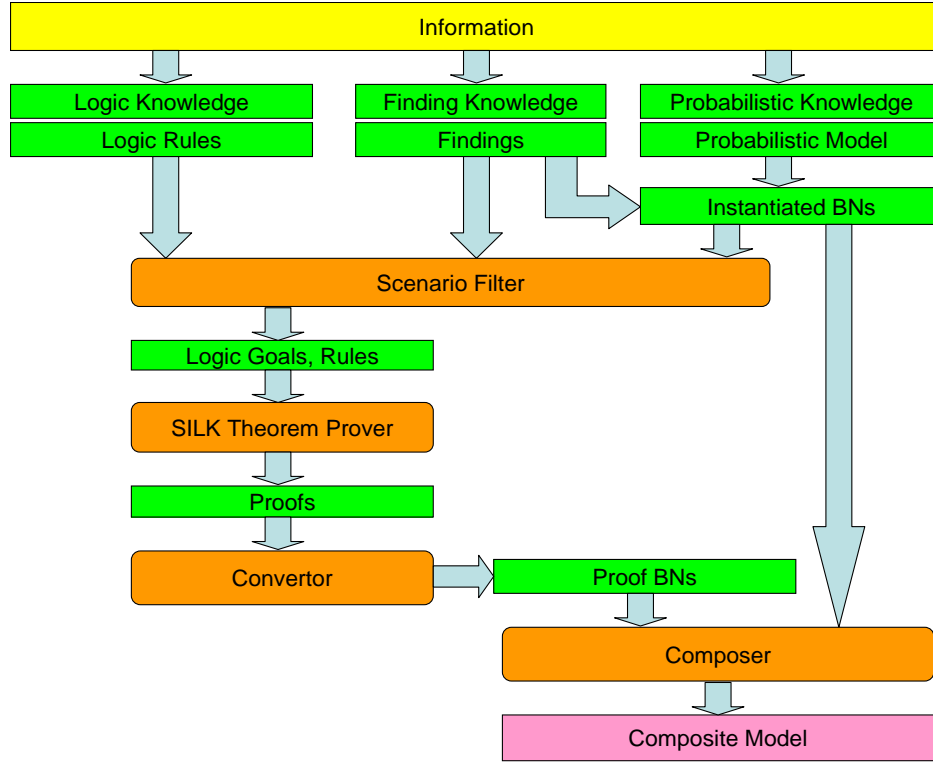


Figure 15. The BALER software process flow, which is supported by the tripartite ontology of real world concepts, events, and information sources

## 2.5 Causality

Causality is a special relationship among events for which certain properties hold probabilistically. For example, causality is logically irreflexive and asymmetric, but probabilistically transitive. Causality, like the relation *subevents*, generates a strict partial order among events. Causal models are very useful, because they allow prediction of the effect of interventions [3,5]. Our interest is in a causal Bayesian network.

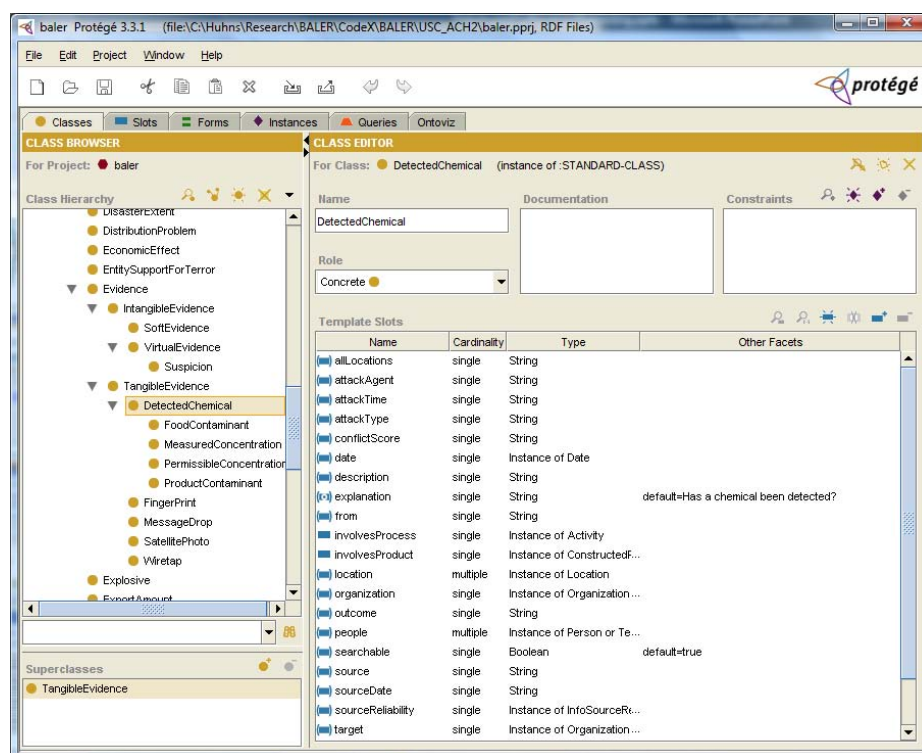
A *causal Bayesian network* consists of a causal graph, a directed acyclic graph (DAG) expressing causal relationships, and a probability distribution respecting the independence relation encoded by the graph [7]. A link between two nodes in a Bayesian network is often interpreted as a causal link. However, this is not necessarily the case. When each link in a Bayesian network is causal, then the Bayesian network is called a causal Bayesian network or Markovian model. A Markovian model is a popular graphical model for encoding distributional and causal relationships. To summarize, a Markovian model consists of a DAG  $G$  over a set of variables  $V = \{V_1; \dots; V_n\}$ , called a causal graph and a probability distribution over  $V$  that has some constraints on it. The interpretation of such a model consists of two parts: the association of the variables to events and the assignment of probability distributions to the links. For causality, variable assignment must satisfy the obvious constraint that

$$(\text{Event A causes Event B}) \rightarrow (\text{time}_A < \text{time}_B)$$

The probability distributions must satisfy two constraints. The first constraint is that each variable in the graph is independent of all its non-descendants given its direct parents. The second constraint is that the directed edges in  $G$  represent causal influences between the corresponding variables. A Markovian model for which only the first constraint holds is called a Bayesian network, and its DAG is called a Bayesian network structure. This explains why Markovian models are also called causal Bayesian networks. As far as the second condition is concerned, causality requires that, when a variable is set, the parents of that variable be disconnected from it: this is called the excision model of causality.

In our prototype tool, Magellan, new variables are added to the causal and event portion of an analyst's ontology using Protégé, so that all of the nodes in a Bayesian network fragment are represented in a standard and consistent terminology. We extend SUMO with this terminology, so that we can take advantage of SUMO's existing description of general knowledge of the world. Each variable has a set of identifying attributes, which are used to combine fragments (fragments can be combined only if their attributes unify) [1,4]. See Figure 16.

Probabilities are assigned to events in the fragment by performing experiments, estimating beliefs, or counting outcomes. Once assigned, they are updated by conditioning on evidence using Bayes rule and the laws of probability. The fragments are stored in a repository, where they can be matched with evidence and combined with other fragments to produce models of situations that are as complete, accurate, and specific as possible.



**Figure 16. Protégé is used to enter the ontology concepts that form the basis for representing situations and evidence**

## 2.6 Evidence

Fragments are instantiated by evidence, which we define informally as information (perhaps wrong, perhaps incomplete) about what happened (events). For example, a bank clerk might be uncertain whether a money transfer was to a Cairo bank or a Boston bank. We represent in the information source ontology the level of credibility of items of evidence, and provide a Bayesian interpretation of credibility. Formally, we define *evidence* to be a collection of findings, each of which describes the state of a Bayesian network variable, and distinguish three kinds [8]:

1. A *hard finding* specifies that the variable has a particular value. For example, whether or not a money transfer occurred or whether or not a suspect is a terrorist

$(Male\_TerroristSuspect = true)$

2. A *soft finding* is a distribution on the states of a variable, usually corresponding to an “objective” statistical distribution that is not expected to change within a scenario. For example, there might be an observation that 95% of terrorists are male (and 5% are not), i.e.,

$Q(Male\_TerroristSuspect) = (0.95, 0.05)$

3. A *virtual finding* is a likelihood ratio corresponding to the credibility associated to an evidence source, such as a witness. For example, witness Bill might have observed a suspect entering a men's-only area of a mosque, which would be interpreted as 4-to-1 that the suspect is a male

$L(Male\_TerroristSuspect) = (0.8, 0.2)$

Unlike soft findings, virtual findings allow for an update of the posterior probability of the evidence variable.

The relationships among the evidence types are shown in Figure 17.

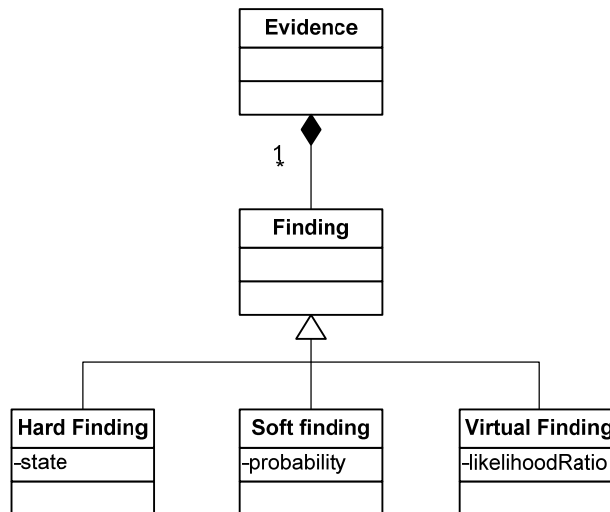
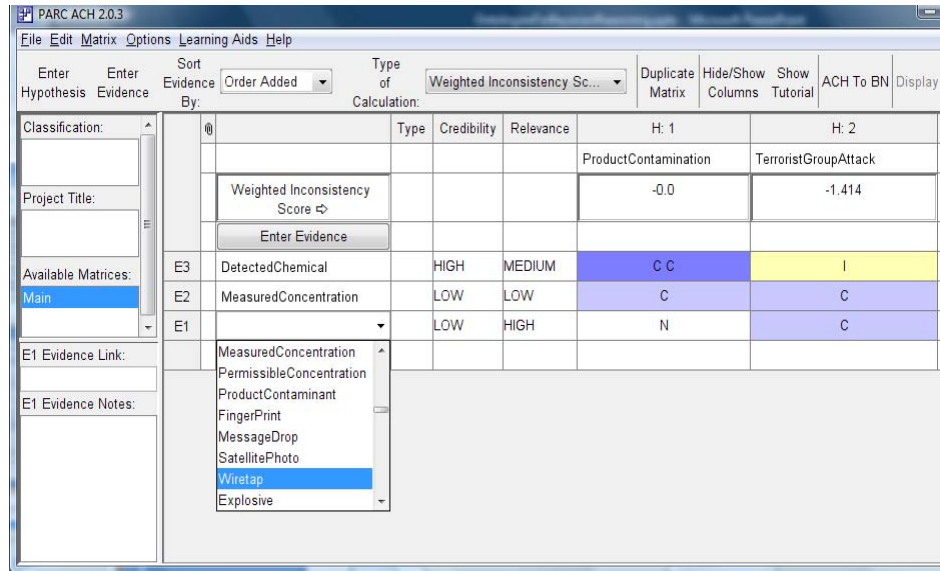


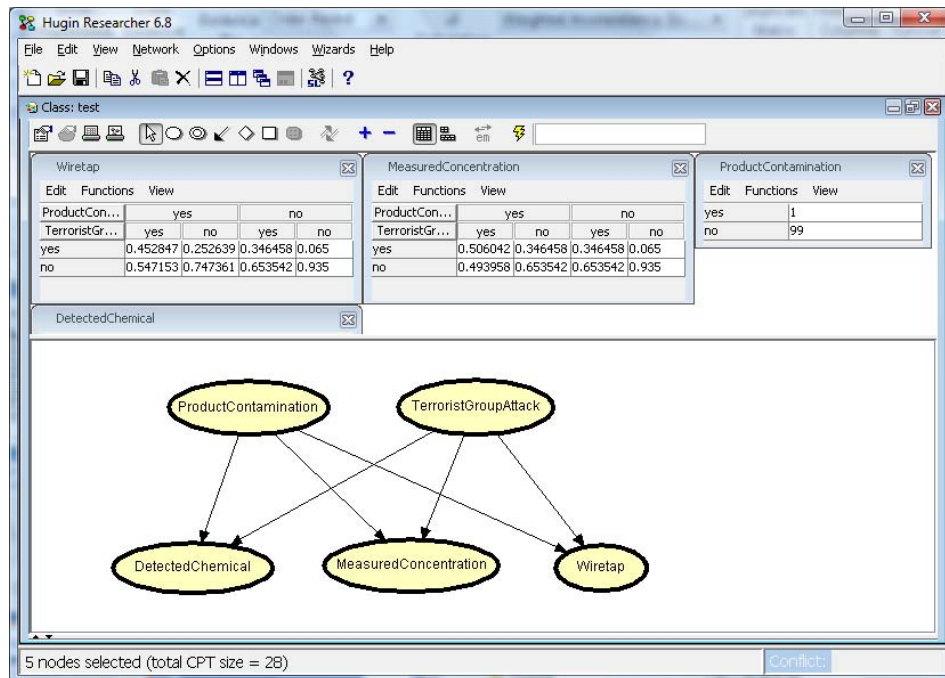
Figure 17. Evidence consists of a set of findings, which can be of three different types



**Figure 18. Magellan’s extended ACH interface is integrated with the ontology of events through pull-down menus**

Our modified version of the tool ACH2 [6] is used by an analyst to enter the appropriate hypotheses and any initial evidence that might be available. The terminology available to the analyst is provided via drop-down menus as shown in Figure 18, where the menu entries are the ontology terms from our ontology developed in Protégé. The resultant Analysis of Competing Hypotheses (ACH) [2] matrix is converted automatically into a bipartite Bayesian network, with initial probabilities assigned based on the relevance factors assigned to cells of the matrix. An example of the network is shown in Figure 19. The network is saved into a repository of fragments, from where it can be retrieved for matching to evidence.





**Figure 19. A Bayesian network fragment constructed automatically from an ACH matrix. The conditional probabilities needed for Bayesian reasoning are derived from the user-entered values in the matrix indicating whether or not a finding is consistent with an analyst's hypothesis.**

### 3.0 Use of Tripartite Ontology for Intelligence Analysis

Figure 20 shows an end-to-end architecture for Bayesian reasoning, which would be used as follows. The process might be triggered by the arrival of evidence in the form of a message, such as the following:

**FBI Report Date: 10 April 2003.** FBI: Abdul Ramazi is the owner of the Select Gourmet Foods shop in Springfield Mall, Springfield, VA. (Phone number 703-659-2317). First Union National Bank lists Select Gourmet Foods as holding account number 1070173749003. Six checks totaling \$35,000 have been deposited in this account in the past four months and are recorded as having been drawn on accounts at the Pyramid Bank of Cairo, Egypt and the Central Bank of Dubai, United Arab Emirates. Both of these banks have just been listed as possible conduits in money laundering schemes.

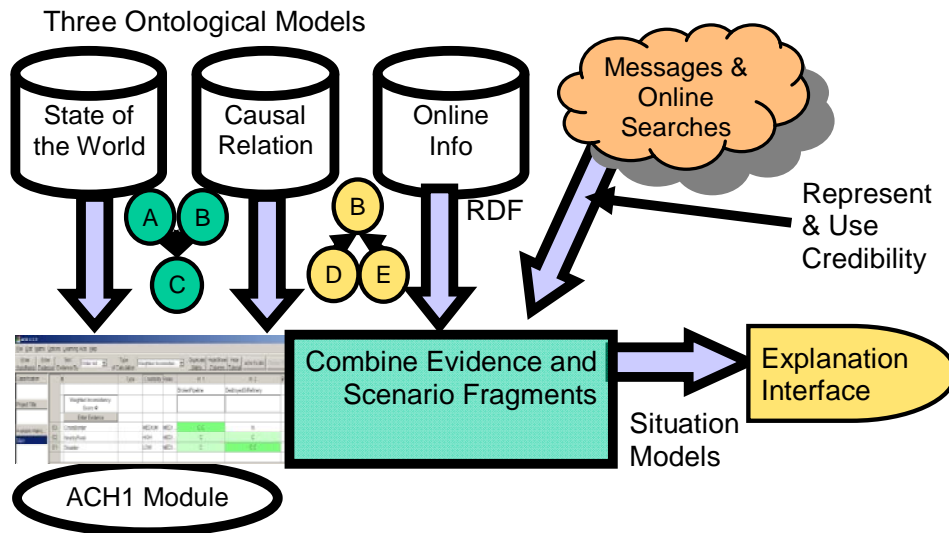


Figure 20. Magellan architecture for Bayesian Reasoning used to explore an analyst's hypotheses

Based on such a message, or based on a hypothesized situation that an analyst would like to investigate, an appropriate scenario represented as a Bayesian model is chosen by the analyst and a corresponding form is displayed listing initial evidence and the domain variables for the scenario. The evidence values for the variables can be supplied automatically from the triggering messages, by matching message terms with ontology concepts as shown in Figure 21, or can be entered by the analyst. Because the probabilities of the variables represented in a situation are updated to be consistent with the evidence at hand, the situation tracks the variables of interest to an analyst. When the probability of a particular value of a variable of interest becomes sufficiently high, an alert could be issued to the analyst.



### Algorithm 1. Value-of-Information Calculation

- Let  $V$  be a variable whose value affects the actions to be taken by an analyst. For example,  $V$  indicates whether a bomb is placed on a particular airliner.
- Let  $p(v)$  be the probability that variable  $V$  has value  $v$ .
- The entropy of  $V$  is:

$$H(V) = - \sum_{v \in V} p(V = v) \log(p(V = v))$$

- Let  $T$  be a variable whose value we may acquire (by expending resources). For example,  $T$  indicates whether a passenger is a known terrorist.
- The entropy of  $V$  given that  $T$  has value  $t$  is:

$$H(V|t) = - \sum_{v \in V} p(V = v|T = t) \log(p(V = v|T = t))$$

- The expected entropy of  $V$  given  $T$  is:

$$E[H(V|t)] = \sum_{t \in T} p(T = t) H(V|t)$$

- The value of information is then:

$$VOI(V) = -(E[H(V|t)] - H(V))$$

Figure 22. The Value of Information calculation algorithm used in Magellan.

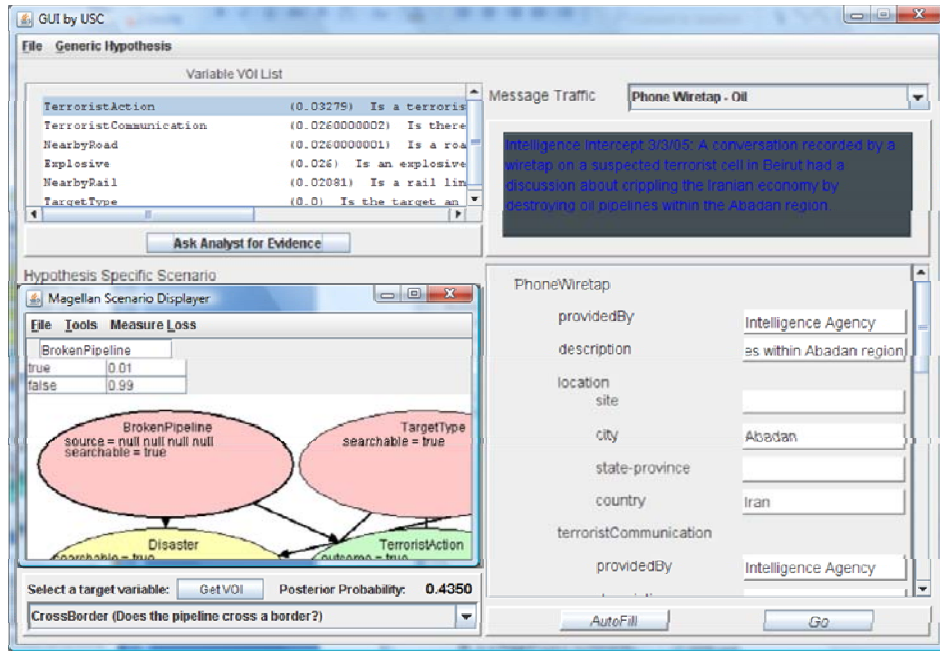


Figure 23. The Magellan interface showing an evidence message, the ontology concepts it contains, the fragments that it instantiates, composed into a situation, and the posterior probability for an hypothesis about the situation.

## 4.0 Evaluation

An early anecdotal evaluation of Magellan was conducted at NIST. The evaluators (three naval reservists with a background in intelligence analysis) tested the hypothesis generation aspect of the system for four hours. In this test, the analysts were presented with several items of evidence (similar to the FBI Report in of section 3) and asked to generate hypotheses, using an interface such as is shown in Figure 22. After they had finished, they were shown hypotheses generated by Magellan and were asked to rate these hypotheses in comparison to the ones they had generated. The NIST summary of the evaluation indicated that the analysts generated more hypotheses than Magellan and that Magellan's hypotheses did not take into account all the possible variables. However, analysts' ratings for Magellan-generated hypotheses are equal to the ratings for the analyst-generated hypotheses in 1/3 of the cases. In 7/9 cases the ratings for the Magellan-generated hypotheses were given mid-level ratings or higher.

## 5.0 Conclusion

Our work is predicated on the observation that ontologies make it easier for tools to interoperate. We have found that our ontologies need to describe both the physical world and the on-line information world, because our reasoning system relies on the relationships and links between both kinds of domains. The reasoner, BALER, enables first-order logic sentences to be combined with Bayesian networks by generating Bayesian networks for any first-order natural deduction proof (that uses the Reeves-Clarke inference rules). This exploits the complementary powers of both logical and Bayesian representations.

Lessons learned from this project

- Causality is a complicated relation
- The probability and logic community is varied and does not accept that Bayesian approaches to probability are universally valid; probability intervals and possibilistic approaches are valued
- It is still an open question whether there is a need to support soft evidence in applications: virtual evidence may be enough
- We did not explore the use of soft evidence in models for IA
- Our approach to translate logic into Bayesian networks is promising, and seems more principled than approaches that add probability intervals to description logics
- The choice of which proofs to translate into Bayesian networks is very difficult, and cannot be avoided without introducing enormous computational complexity
- Composition is harder computationally than expected--better heuristics are needed for focusing search.

### 5.1.1 The Relationship between Proofs and Bayesian Networks

Our proposal for CASE was motivated by the observations that:

- Many problems contain both logical and probabilistic aspects
- Proofs and Bayesian Networks (BNs) are both directed acyclic graphs, sometimes over logical formulas, and therefore should be combinable.

The exact nature of the relationship became clear only after some investigation. The intuitive direction for edges in proofs is from premises to conclusions, but this direction violates independence relationships for BNs. The direction that properly captures independence is to direct edges from subformulas to superformulas. The conditional probability tables (CPTs) for the internal nodes are then the simple binary-valued logical connective definitions (in the propositional case).

### 5.1.2 Proofs and Theories

When the superformulas are fully expanded, the proof structure seems to disappear and we seem to have constructed BNs from theories rather than from proofs, in which case the proof search seems superfluous. This was observed by peers when we tried to present our work, and we need to present examples and use cases which help to clarify the differences.

- Firstly, any proof is simply a set of subformulas and superformulas; this is brought out most explicitly by normal natural deduction proofs as we used in the project.
- Secondly, the axioms of the theory are not necessarily the largest superformulas; the superformulas may need to be synthesized as part of the argument. This is the explicit structure of typical geometry proofs (in which one first builds additional structure beyond what was given in the premises so that the total structure may then be analyzed in order to deduce properties of the original given structure). When we choose overly simple example problems (as could be handled by prolog, for example) then we fail to highlight the role of synthesis in the reasoning.
- The next point to recognize is that in the case of very large theories, the proofs pull out small parts (sub theories) which are sufficient to support the reasoning. This point becomes much more significant in first-order theories, as discussed below.

### 5.1.3 First-order Theories

We started work with propositional proofs because these are easier to think about and represent. Of course, BNs handle propositional variables fairly directly, so the real gains come when dealing with first-order theories. There is a wide variety of approaches to extending propositional variables to have attributes, but none seem to capture full first order logic (especially nested quantifiers). Our work can be differentiated primarily by its expressivity: we do capture full first-order logic. This part of the theory needs to be verified (I have a sketch of a proof that was never finalized in full generality). Once verified, this feature should be stated up front, as some of our peers have had difficulty in understanding some of our goals. The main issue raised in the literature is complexity of updating the networks. Poole 2003, “First-Order Probabilistic Inference”, and de Salvo Braz et al 2007, “Lifted First-Order Probabilistic Inference” are probably the most relevant to our work, and discuss complexity explicitly. We should compare complexity. We should also look more into Markov Logic Networks in order to gain a deeper understanding of the relations and differences between our approach and MLNs.

### 5.1.4 Semantics

In order to use traditional Tarski models for formal logic, we interpret  $P(A) = 90\%$  as referring to a probability distribution over sampling from possible worlds in which the probability a drawing a world in which  $A$  is true is 90%. I believe this comes from Halpern 1990, “An analysis of first-order logics of probability”, which we should study and cite. When  $A$  is the sentence, “Tweety is a bird”, it is hard to imagine what other interpretation might be useful.

However, when  $A$  is the sentence “All birds fly”, then an alternative interpretation is one in which possible worlds come with probability distributions, and each satisfying world has a distribution  $P$  such that  $P(\text{fly} \mid \text{bird}) = 90\%$ . This second interpretation is discussed in Bacchus 1990, *Representing and Reasoning with Uncertain Knowledge*.

The first is certainly easier to work with and is fully general over first-order logic formulas, but the second might better capture our intuitions in certain situations. Under the first interpretation, we are 90% certain that absolutely every bird flies, but in the 10% likely case that not every bird flies, there is no distinction at all between worlds in which no birds fly and worlds in which all birds but one fly.



## 6.0 References

- [1] John Cheng, Ray Emami, Larry Kerschberg, Eugene Santos, Jr., Qunhua Zhao, Hien Nguyen, Hua Wang, Michael Huhns, Marco Valtorta, Jiangbo Dang, Hrishikesh Goradia, Jingshan Huang, and Sharon Xi, “OmniSeer: A Cognitive Framework for User Modeling, Reuse of Prior and Tacit Knowledge, and Collaborative Knowledge Services,” *Proceedings of the 38th Hawaii International Conference on System Sciences* HICSS38, 2005.
- [2] Richards J. Heuer, Jr., *Psychology of Intelligence Analysis*, Center for the Study of Intelligence (at <http://www.cia.gov/csi/books/19104/index.html>), 1999.
- [3] Yimin Huang and Marco Valtorta, Pearl’s Calculus of Intervention is Complete, *Proceedings of the Twenty-Second Conference on Uncertainty in Artificial Intelligence (UAI-06)* (2006), 217-224.
- [4] Katherine Laskey and Suzanne Mahoney, Network Fragments: Representing Knowledge for Constructing Probabilistic Models, *Proceeding of the Thirteenth Conference on Uncertainty in Artificial Intelligence* (1997), 334-341.
- [5] Judea Pearl, *Causality: Modeling, Reasoning, and Inference*, Cambridge, England: Cambridge University Press, 2000.
- [6] Peter Pirolli and Lance Good, Evaluation of a Computer Support Tool for Analysis of Competing Hypotheses, UIR Technical Report, Palo Alto Research Center (), 2004.
- [7] Marco Valtorta, John Byrnes, and Michael Huhns, Logical and Probabilistic Reasoning to Support Information Analysis in Uncertain Domains, *Proceedings of the Third Workshop on Combining Probability and Logic (Prolog-07)*, Canterbury, England, September 5-7, 2007
- [8] Marco Valtorta and Yimin Huang, Identifiability in Causal Bayesian Networks: A Gentle Introduction, *Cybernetics and Systems* 39:4 (2008), 425-442.
- [9] Marco G. Valtorta, Y.-G. Kim, and Jirka Vomlel, Soft Evidential Update for Multiagent Systems, *International Journal of Approximate Reasoning* 29:1 (2002), 71-106.

## 7.0 List of Acronyms

ACH.....	Analysis of Competing Hypotheses
AFRL.....	Air Force Research Laboratory
BALER.....	Bayesian and Logical Engine for Reasoning
BC-Hugin.....	Big Clique Hugin
BN.....	Bayesian Network
BRUSE.....	Bayesian Reasoning Using Soft Evidence
CASE.....	Collaboration and Analyst System Effectiveness
CLIPS.....	C Language Integrated Production System
CPT.....	Conditional Probability Table
DAG.....	Directed Acyclic Graph
FBI.....	Federal Bureau of Investigation
IA.....	Information Assurance
JESS.....	Java Expert System Shell
MLN.....	Markov Logic Network
NIST.....	National Institute of Standards and Technology
RDF.....	Resource Description Framework
SILK.....	Semantic Inferencing on Large Knowledge
SUMO.....	Suggested Upper Merged Ontology
SPARQL.....	SPARQL Protocol and RDF Query Language
XMLBIF.....	eXtensible Markup Language Bayesian Interchange Format